

1 Patterns of putative gene loss suggest rampant developmental system drift in nematodes

2

3 Gavin C. Woodruff<sup>1\*</sup>

4

5 <sup>1</sup>Institute of Ecology and Evolution, University of Oregon, Eugene, OR, USA

6 \*Correspondence: [gavincw@uoregon.edu](mailto:gavincw@uoregon.edu)

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47 **Abstract**

48  
49 Gene loss often contributes to the evolution of adaptive traits. Conversely, null mutations  
50 frequently reveal no obvious phenotypic consequences. How pervasive is gene loss, what kinds  
51 of genes are dispensable, and what are the consequences of gene loss? The nematode  
52 *Caenorhabditis elegans* has long been at the forefront of genetic research, yet only recently have  
53 genomic resources become available to situate this species in its comparative phylogenetic and  
54 evolutionary context. Here, patterns of gene loss within *Caenorhabditis* are evaluated using 28  
55 nematode genomes (most of them sequenced only in the past few years). Orthologous genes  
56 detected in every species except one were defined as being lost within that species. Putative  
57 functional roles of lost genes were determined using phenotypic information from *C. elegans*  
58 WormBase ontology terms as well as using existing *C. elegans* transcriptomic datasets. All  
59 species have lost multiple genes in a species-specific manner, with a genus-wide average of  
60 several dozen genes per species. Counterintuitively, nearly all species have lost genes that  
61 perform essential functions in *C. elegans* (an average of one third of the genes lost within a  
62 species). Retained genes reveal no differences from lost genes in *C. elegans* transcriptional  
63 abundance across all developmental stages when considering all 28 *Caenorhabditis* genomes.  
64 However, when considering only genomes in the subgeneric *Elegans* group, lost genes tend to  
65 have lower expression than retained genes. Taken together, these results suggest that the genetics  
66 of developmental processes are evolving rapidly despite a highly conserved adult morphology  
67 and cell lineage in this group, a phenomenon known as developmental system drift. These  
68 patterns highlight the importance of the comparative approach in interpreting findings in model  
69 systems genetics.

70

71

72 **Keywords**

73

74 Developmental system drift, comparative genomics, gene loss, *Caenorhabditis*

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

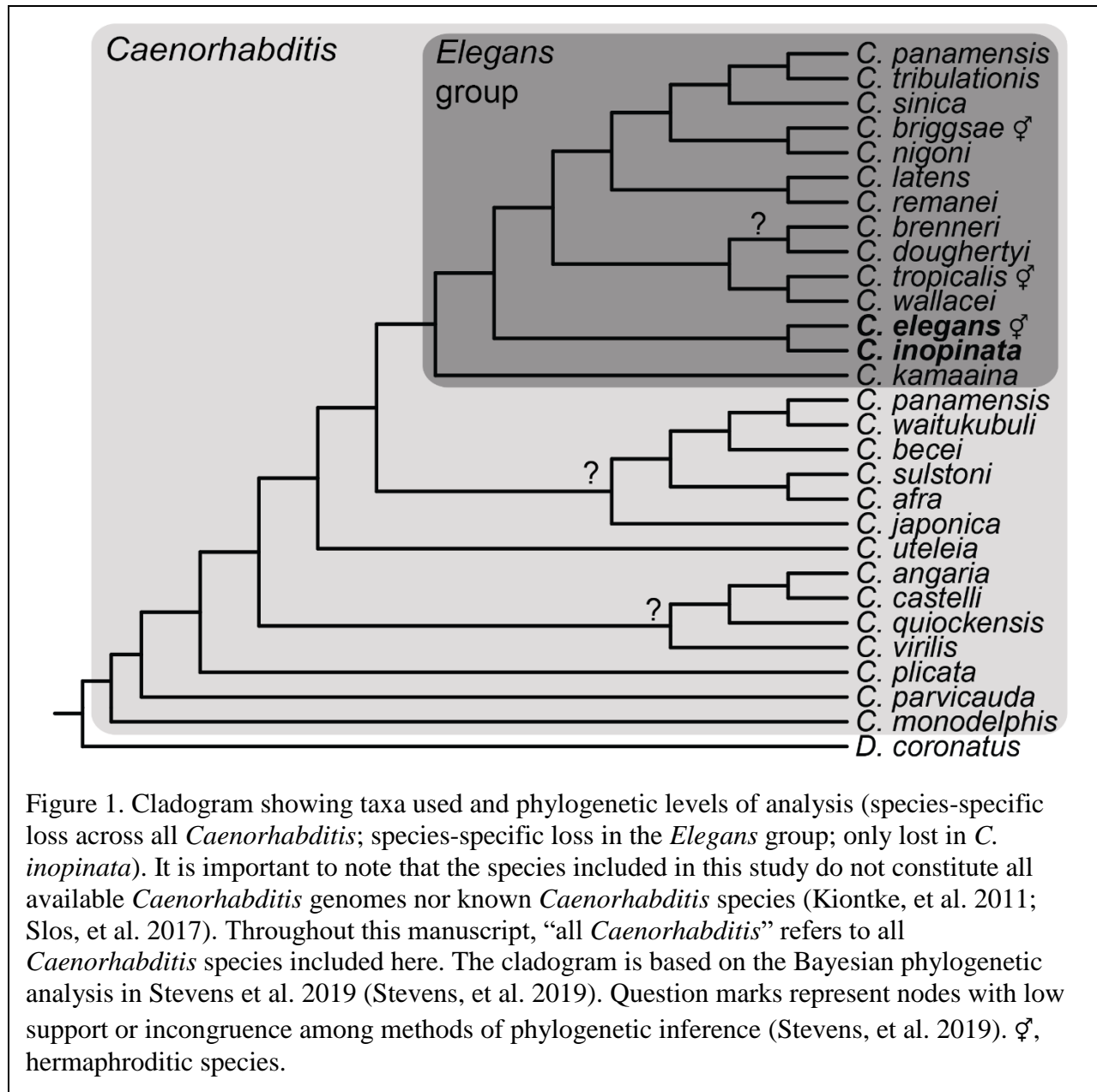
91

## 92 Introduction

93 Gene loss is common and often has phenotypic consequences. Such losses, whether due to large-  
94 scale structural variation or to single nucleotide changes that render proteins non-functional, are  
95 typically associated with disease states that can represent profound public health challenges  
96 (Stankiewicz and Lupski 2010). However, gene loss also frequently underlies adaptive change  
97 (Albalat and Cañestro 2016). Such losses underlie changes in leaf morphology among  
98 Brassicaceae plant species (Vlad, et al. 2014), cold temperature resistance in flies (Greenberg, et  
99 al. 2003), self-incompatibility in *Arabidopsis* (Shimizu, et al. 2008), and pigmentation variation  
100 in multiple systems (Zufall and Rausher 2004; Protas, et al. 2006; Hoballah, et al. 2007).  
101 Similarly, selection can drive gene loss or genome size reduction in the context of experimental  
102 evolution as well (Nilsson, et al. 2005; Good, et al. 2017). The absence of a gene can even  
103 promote reproductive isolation and thereby play an important role in speciation (Bikard, et al.  
104 2009; Ben-David, et al. 2017). Thus, gene loss must contribute to evolutionary change. What  
105 kinds of genes are dispensable, and how might they promote phenotypic divergence?

106 Conversely, although gene loss often has dramatic phenotypic consequences, a common outcome  
107 of gene loss is no observable phenotypic consequence at all. For instance, although the average  
108 human being is homozygous null for about twenty genes, most people do not have genetic  
109 disorders (MacArthur, et al. 2012). In addition, multiple large-scale knockout and knockdown  
110 screens for genetic function have unearthed thousands of genes with no obvious function in  
111 multiple organisms (Winzeler, et al. 1999; Kamath, et al. 2003; Dietzl, et al. 2007). Such  
112 observations are often explained by genetic redundancy, wherein multiple genes perform the  
113 same function, and therefore the loss of any one such gene is of little phenotypic consequence  
114 (Nowak, et al. 1997). However more recent studies have revealed that the fitness consequences  
115 of many gene knockdowns vary depending on genetic background (Dowell, et al. 2010; Chari  
116 and Dworkin 2013; Paaby, et al. 2015). Here such results could be explained by pervasive  
117 compensatory change and developmental system drift (True and Haag 2001), wherein dramatic  
118 differences in underlying developmental processes nonetheless promote similar phenotypes.  
119 Overall, then, the extent to which gene loss influences phenotypic change (or lack thereof) is not  
120 completely understood.

121 The first metazoan to have its genome sequenced, the nematode *C. elegans* has been a widely  
122 used model system for decades (Corsi, et al. 2015). In addition to its widespread use in  
123 developmental and molecular genetics, much is known about its genomic features (Gerstein, et  
124 al. 2010). Indeed, the WormBase database contains vast information for many of its ~20,000  
125 protein-coding genes (Howe, et al. 2016). Despite this, the comparative and evolutionary  
126 genomic resources of *C. elegans* have been historically lacking compared to other widely studied  
127 model systems such as *Drosophila* (Consortium 2007; Huang, et al. 2014; Casillas and  
128 Barbadilla 2017). However, this has recently changed with the rapid discovery of dozens of  
129 *Caenorhabditis* species (Kiontke, et al. 2011; Ferrari, et al. 2017; Slos, et al. 2017; Stevens, et al.  
130 2019) in tandem with the sequencing of multiple close relatives of *C. elegans* (Fierst, et al. 2015;  
131 Slos, et al. 2017; Kanzaki, et al. 2018; Ren, et al. 2018; Rödelsperger 2018; Yin, et al. 2018;  
132 Stevens, et al. 2019). Here, I combine the collective knowledge of *C. elegans* developmental  
133 genetics that resides in the WormBase database with patterns of gene loss observed across the  
134 genomic evolution of 28 *Caenorhabditis* species, finding that patterns of species-specific gene



135 loss underlie vast developmental system drift in this genus. These patterns underscore the crucial  
136 role of genomic context in understanding gene function.

## 137 **Results**

138 *All Caenorhabditis species have lost multiple genes that perform essential functions in C. elegans*

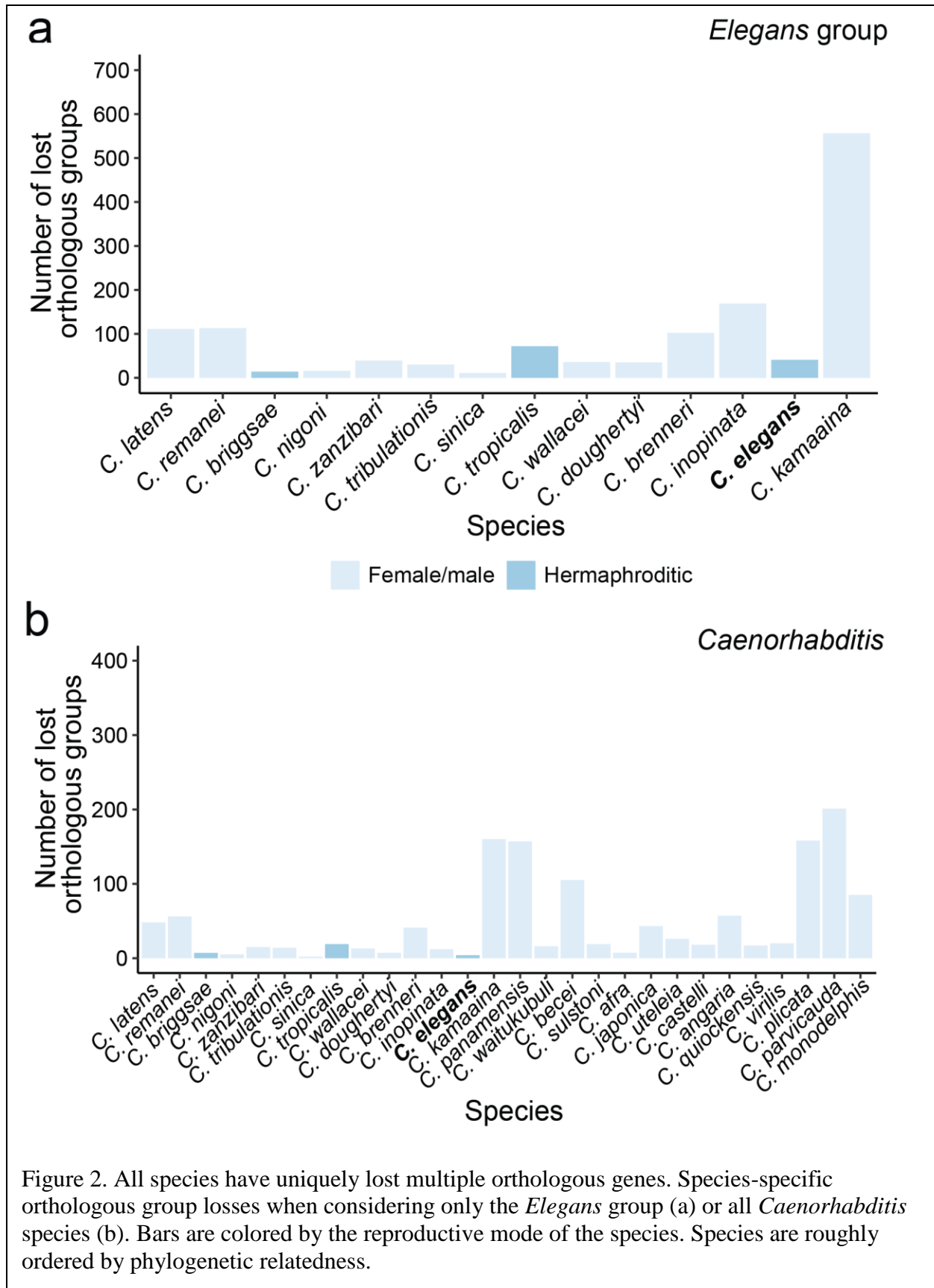
139 To explore the extent and consequences of gene loss in *Caenorhabditis*, species-specific gene  
140 losses were inferred at two levels of phylogenetic scope (the whole genus and only the *Elegans*  
141 group, Figure 1). Briefly, after defining groups of orthologous proteins across 28 *Caenorhabditis*  
142 species (Emms and Kelly 2015), orthologous groups present in all species but one were called  
143 presumptive species-specific lost genes. Here, gene loss will be assumed to be equivalent to this

144 type of species-specific absence, as opposed to many other possible patterns of repeated loss in  
145 multiple species; here I only examined patterns of species-specific loss.

146 At both phylogenetic levels, all species exhibit multiple species-specific gene losses (Figure 2).  
147 As the number of shared orthologous groups declines as more species are included  
148 (Supplemental Figure 6), the number of gene losses per species is higher when considering the  
149 *Elegans* group (mean=96, median=40, range=11-556) than when considering the genus as a  
150 whole (mean=48, median=19, range=2-201). As the genome assemblies under consideration are  
151 at varying degrees of completeness and quality (Supplemental Figures 3-5), this may have  
152 influenced the number of inferred species-specific losses. However, gene loss is only  
153 significantly associated with genome completeness when considering the *Elegans* group  
154 (Supplemental Figure 8) and not the whole genus (Supplemental Figure 9). Furthermore, gene  
155 loss is not significantly correlated with scaffold number (Supplemental Figures 10-11) or N50  
156 (Supplemental Figure 12-13). Thus, although genome quality may influence the inference of  
157 gene loss, and the results here may overestimate gene loss, most genome quality metrics are not  
158 correlated with gene loss. Moreover, there are a number of species with high quality reference  
159 genomes (*C. elegans*, *C. briggsae*, *C. tropicalis*, *C. nigoni*, *C. wallacei*, and *C. inopinata*), and  
160 all of these species exhibit species-specific gene loss (Figure 2). Indeed, in the case of *C.*  
161 *inopinata*, the degree of gene loss with respect to the *Elegans* group is high (169 lost orthologous  
162 groups), despite its completely assembled reference genome (seven, chromosome-level  
163 scaffolds) and high BUSCO completeness score (98.1%)(Kanzaki, et al. 2018). So although  
164 genome quality likely inflates the extent of gene loss in some species, patterns of species-specific  
165 gene loss are still detected even in high quality assemblies.

166 To understand their functional relevance, lost genes were paired with WormBase phenotype data  
167 ((Schindelman, et al. 2011); see methods). WormBase is a repository for biological knowledge of  
168 *C. elegans*, notable for housing various kinds of genomic data related to *C. elegans* and its  
169 relatives (Howe, et al. 2015). Among these are “phenotype” terms, which constitute a formal  
170 ontology used to describe phenotypes associated with genes (Schindelman, et al. 2011). More  
171 specifically, these describe biological phenotypes that arise upon some perturbation of a given  
172 gene, usually through mutation or RNAi knockdown (Schindelman, et al. 2011). There are at  
173 least 2,443 phenotype terms in WormBase, ranging from the straightforward (“embryonic  
174 lethal,” “no germ line”) to the esoteric (“loss of asymmetry AWC,” “nuclear fallout”). I paired  
175 all of the *C. elegans* gene constituents of lost orthologous groups in each species at both  
176 phylogenetic levels with their WormBase phenotype terms. Among genes with phenotypes (only  
177 about 42% of *C. elegans* protein-coding genes were found to have phenotypes in WormBase  
178 (8,514 out of 20,204)), I further classified them into two categories: essential and inessential.  
179 Essential phenotypes were defined by the presence of any of these words: “lethal,” “arrest,”

180



181 “sterile,” or “dead,” and ultimately constituted 58 unique phenotype terms (see Supplemental  
182 Data for list of essential phenotypes). All other phenotypes were noted as inessential. It is  
183 important to note that there are multiple phenotypes here noted as inessential that are probably  
184 critical for survival and that these numbers of essential genes reported here are likely  
185 underestimates. Additionally, the phenotypes of genes lost only in *C. elegans* cannot be assessed  
186 because WormBase phenotypes are derived from studies of genes that are present in *C. elegans*.

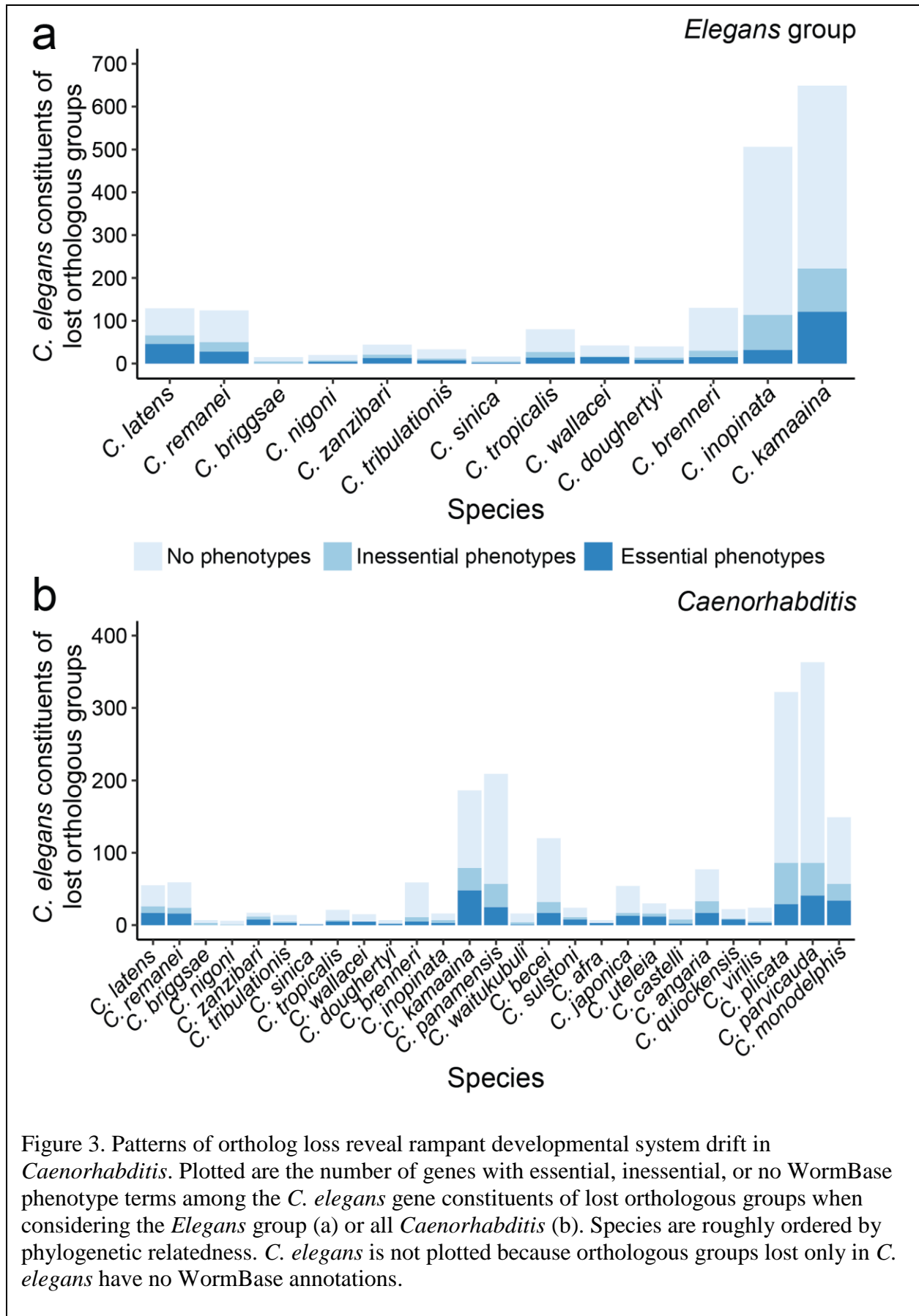
187 At both phylogenetic levels considered, nearly all species have lost genes that are associated with  
188 both essential and inessential phenotypes (Figure 3). Among the *Elegans* group, 36% of the *C.*  
189 *elegans* genes associated with species-specific lost orthologous groups had phenotypes on  
190 average (range=23%-53%), and 20% had essential phenotypes (range=0%-36%; Figure 3a).  
191 Across all *Caenorhabditis*, 37% of the *C. elegans* genes associated with species-specific lost  
192 orthologous groups had phenotypes on average (range=17%-71%), and 23% had essential  
193 phenotypes (range=0%-50%; Figure 3b). Notably, *C. briggsae* has not lost any essential genes at  
194 both levels of phylogenetic consideration; its close relative *C. nigoni* has also lost no essential  
195 genes when considering the whole genus (although it has lost four essential genes when  
196 considering the *Elegans* group; Figure 3). All other species have lost genes that are needed for  
197 viability and fecundity in *C. elegans*. These patterns suggest that genetic functions among highly  
198 conserved processes (such as embryogenesis) are rapidly evolving in this group.

#### 199 *Patterns of pleiotropy, specificity, and spatiotemporal transcript abundance among lost genes*

200 Do lost genes share any common features, or can gene loss be predicted? To address this  
201 question, other features of *C. elegans* genes associated with species-specific lost orthologous  
202 groups were also examined. In addition to phenotypes, WormBase also contains “anatomy” and  
203 “life stage” ontologies. These relate spatial (“anatomy”) and temporal (“life stage”) expression  
204 patterns to genes. WormBase also contains information about pairwise interactions among genes  
205 (“interaction”), which are defined by epistatic genetic interactions or physical/biochemical  
206 interactions. WormBase also tracks the number of peer-reviewed scientific papers that mention a  
207 given gene as a “reference count.” Additionally, the domains in all *C. elegans* proteins were  
208 defined using the 16,713 Pfam domain seed alignments and HMMER (Finn, et al. 2015). The  
209 per-gene number of unique features of all of these categories (phenotype, tissue, life stage,  
210 interaction, reference count, and domain) were counted. This then provides quantitative measures  
211 of: the consequential phenotypic complexity upon perturbation of a given gene (phenotype); the  
212 extent of expression specificity across space and time of a given gene (anatomy and life stage);  
213 the connectedness of a given gene in its biological network (interaction); the extent to which a  
214 given gene has been studied (reference count); and the number of functional modules a given  
215 gene has (domain). Taken together, these provide various coarse measures of genetic specificity  
216 and pleiotropy.

217 All *C. elegans* genes associated with species-specific lost orthologous groups, irrespective of  
218 species, were denoted as lost at the two levels of phylogenetic scope (*Elegans* group or all  
219 *Caenorhabditis*). In addition, genes lost only in *C. inopinata* (within the context of the *Elegans*  
220 group; Figure 2a and Figure 3a) were also addressed. *C. inopinata* is a species worthy of  
221 consideration on its own for two major reasons. First, it is the closest known relative of *C.*  
222 *elegans* (Kanzaki, et al. 2018; Woodruff, et al. 2018) and represents the lower bound of  
223 phylogenetic distance from the reference species among the organisms in this study. Second, it is  
224 morphologically and ecologically divergent from *C. elegans*, and genes lost only in *C. inopinata*  
225 may be good candidates for understanding the genetic basis of morphological and ecological

226





227 divergence (Kanzaki, et al. 2018; Woodruff and Phillips 2018; Woodruff, et al. 2018). In any  
228 case, at all levels of phylogenetic consideration, genes that were not denoted as “lost” were  
229 called “retained.” Then, the distributions of numbers of WormBase phenotypes, anatomies, life  
230 stages, interactions, reference counts, and PFAM domains among lost and retained genes at the  
231 three levels of phylogenetic scope were compared.

232 The total number of *C. elegans* genes from lost orthologous groups is substantial when  
233 considering both the *Elegans* group (1,828 or 9.0% of *C. elegans* protein-coding genes) and all  
234 *Caenorhabditis* (1,903 or 9.4%). Furthermore, although these genes represent similar proportions  
235 of the genome, they largely do not overlap (only 464 genes are shared among the two groups  
236 (464/3,267 or 14.2%); Supplemental Figure 16). In the case of only *C. inopinata*, the number of  
237 lost *C. elegans* genes is far less (506 or 2.5%). The distribution of ontological term numbers  
238 among lost and retained genes described above also varies depending on the phylogenetic scope  
239 (Figure 4a). Among genes lost only in *C. inopinata*, the average number of all WormBase terms  
240 and domains are significantly lower than among those genes that are retained (Figure 4a). When  
241 considering genes lost among *Elegans* group members, this pattern is similar, although the effect  
242 size of gene loss is far less across all categories (Figure 4a). However, when looking at the  
243 broadest phylogenetic level, all *Caenorhabditis*, this pattern is largely eroded. Lost genes at this  
244 level are largely no different from retained genes; however, lost genes reveal a subtle but  
245 detectable increase in the number of domains relative to retained genes (Cohen’s *d* effect  
246 size=0.073 ; 95% confidence interval=0.019-0.13; Mann-Whitney U  $p=2.09 \times 10^{-10}$ ,  
247  $W=11584369$ ). Thus, the degree of specificity and pleiotropy among lost genes depends upon the  
248 phylogenetic context considered. At narrower phylogenetic scopes, lost genes tend to be less  
249 pleiotropic and specific than retained genes, and as the phylogenetic scope broadens, lost genes  
250 tend to resemble retained genes.

251 In addition to the ontological information accessible in WormBase, the transcriptome of *C.*  
252 *elegans* has also been intensively studied (Gerstein, et al. 2010; Levin, et al. 2012; Hashimshony,  
253 et al. 2015). One recent report measured transcript levels at 30-minute intervals across  
254 embryonic development to define the “time-resolved transcriptome of *C. elegans*” (Boeck, et al.  
255 2016). In addition to patterns of transcription across embryogenesis, this study also included  
256 measures of gene expression across various postembryonic stages (Boeck, et al. 2016). To  
257 provide further biological context to the lost genes defined above, the transcriptomic data from  
258 the Boeck et al. study was paired with this information (Figures 4b-c). Much like with the  
259 WormBase ontological terms (Figure 4), the transcriptional abundances of lost and retained  
260 genes at various embryonic (Figure 4b) and postembryonic (Figure 4c) stages were compared  
261 within the context of three phylogenetic levels (only *C. inopinata*, the *Elegans* group, and all  
262 *Caenorhabditis* (Figure 4b-c)). Like the WormBase ontological terms, patterns of gene  
263 expression among lost genes varied depending on the phylogenetic scope. Among genes lost  
264 only in *C. inopinata*, lost genes exhibited much lower expression than retained genes across all  
265 developmental stages (average effect size= -0.72; Figure 4b-c). Among genes lost in *Elegans*  
266 group species, lost genes are only slightly less expressed than retained genes at all developmental  
267 stages (average effect size= -0.11; Figure 4b-c). And, as with the WormBase terms, no  
268 differences in expression among lost and retained genes could be detected at any developmental  
269 stage at the broadest phylogenetic scope (all *Caenorhabditis*; Figure 4b-c). Thus, transcriptional  
270 abundance among genes with the capacity to be lost also depends on phylogenetic context, and at  
271 narrower phylogenetic scopes, lost genes tend to have lower expression than retained genes.

272

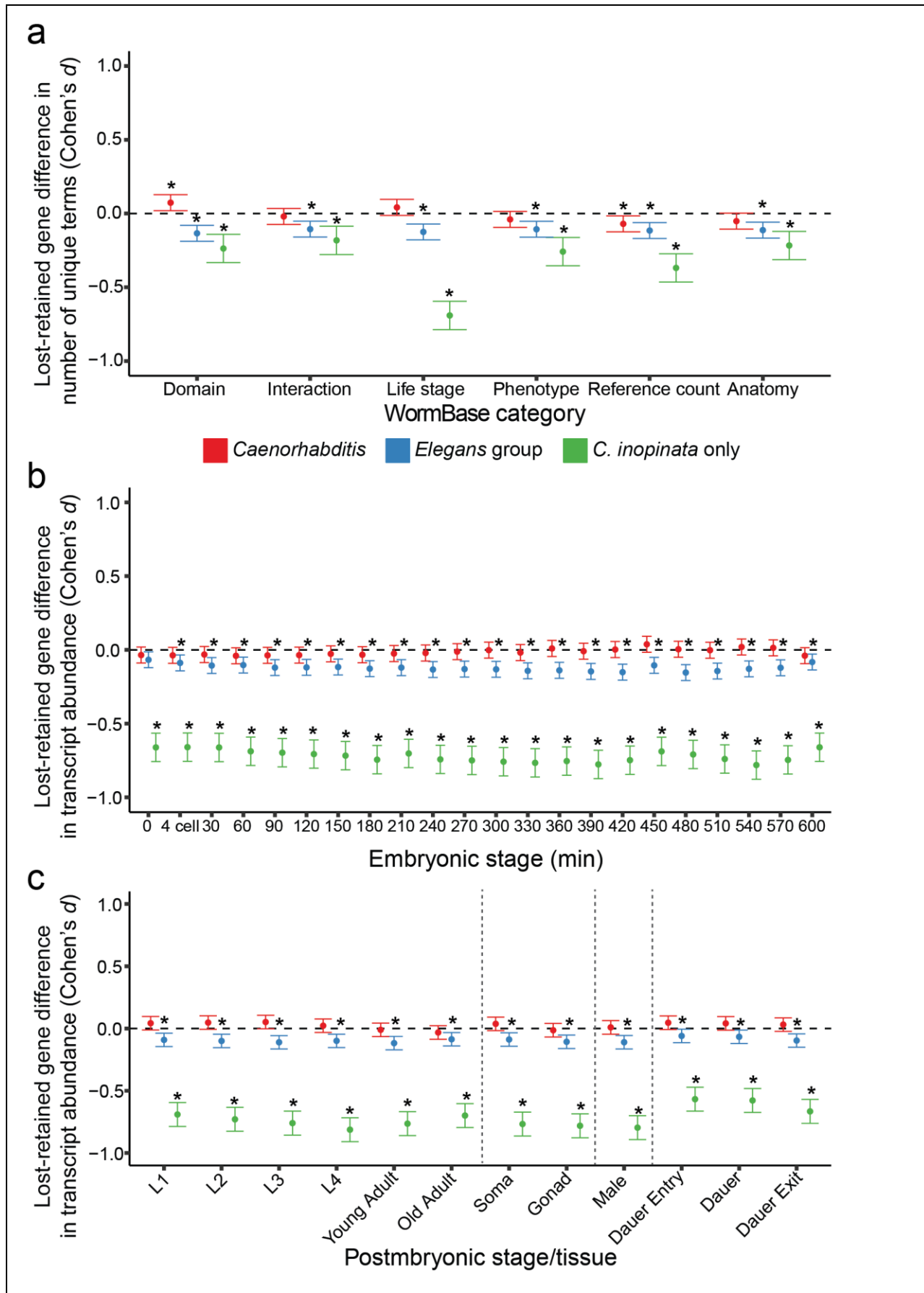


Figure 4. The impact of pleiotropy and transcription on gene loss depends on phylogenetic scope. In all panels the effect size (Cohen's  $d$ ) of gene loss relative to gene retention is plotted. Here, all species-specific gene losses are pooled and denoted as lost genes; *C. elegans* measures among lost and not lost (i.e., retained) genes are being compared. An effect size of -1 notes that the average value among lost genes is one pooled standard deviation lower than retained genes; an effect size of 0 (dashed horizontal line) reveals on average no difference between lost and retained genes. Error bars represent 95% confidence intervals of the effect size. Distributions of all categories for lost and retained genes across all levels of phylogenetic scope can be found in Supplemental Figures 17-25. In all panels, \* = Mann-Whitney U  $p < 0.01$  in a comparison of lost and retained genes. (a) The effect size of gene loss on WormBase feature number per gene. Numbers of domains were determined with HMMER. All other features were retrieved from WormBase. (b) The effect size of gene loss on transcript abundance ( $1 + \log_2(\text{dcpm})$ ) per gene across stages of embryonic development. All stages were recorded as minutes past fertilization with the exception of "4 cell;" the four cell stage is typically ~30 minutes past fertilization in typical rearing conditions (Altun, et al. 2002-2016). (c) The effect size of gene loss on mean transcript abundance ( $1 + \log_2(\text{dcpm})$ ) per gene across stages of postembryonic development. Vertical dotted lines separate hermaphroditic postembryonic stages, hermaphroditic adult soma and germ line preparations, male, and dauer-related stages. RNAseq data were retrieved from (Boeck, et al. 2016).

273 Multivariate analyses were also performed to test the impact of transcriptional activity and the  
274 number of WormBase ontology terms on gene loss. All models of gene loss are significant when  
275 including all transcription and ontology count variables simultaneously (MANOVA: all  
276 *Caenorhabditis*  $p = 5.5 \times 10^{-7}$ , Pillai's trace=0.0050; *Elegans* group  $p = 2.8 \times 10^{-13}$ , Pillai's  
277 trace=0.0072; *C. inopinata* only  $p < 2.2 \times 10^{-16}$ , Pillai's trace=0.021; see supplemental data). The  
278 most important contributors to gene loss (i.e., the factors with the largest coefficients) depend on  
279 the phylogenetic scope. For all *Caenorhabditis*, these are transcription at 330 minutes-past-  
280 fertilization (mpf; Linear Discriminant 1 (LD1) coefficient= -0.49), adult soma-specific  
281 expression (LD1 coefficient= 0.44), and transcription at 300 mpf (LD1 coefficient=0.38; see  
282 supplemental data). For the *Elegans* group, these are young adult hermaphrodite expression  
283 (LD1 coefficient=-0.38), transcription at 420 mpf (LD1 coefficient= -0.36), and transcription at  
284 300 mpf (LD1 coefficient=0.35; see supplemental data). And for *C. inopinata* only, these are  
285 transcription at 480 mpf (LD1 coefficient=0.24), L2 expression (LD1 coefficient=0.22), and  
286 transcription at 390 mpf (LD1 coefficient= -0.22; see supplemental data). However, the  
287 proportion of the variance explained of these models is small (all *Caenorhabditis*, logistic  
288 regression pseudo- $R^2=0.013$ ; *Elegans* group, logistic regression pseudo- $R^2=0.019$ ), although the  
289 models perform better for genes lost only in *C. inopinata* (logistic regression pseudo- $R^2=0.12$ ).  
290 These are consistent with the high overlap of lost and retained genes in principal components  
291 (Supplemental Figures 26-28) and linear discriminant (Supplemental Figure 29) space. Thus  
292 there is a generally weak but detectable impact of transcriptional activity and gene ontology  
293 count on gene loss that increases as the phylogenetic scope narrows.

## 294 Discussion

295 *Widespread turnover of developmental genetic systems in a group with highly conserved*  
296 *morphology*

297 Gene loss is a widespread driver of phenotypic change (Albalat and Cañestro 2016). At the same  
298 time, genetic perturbations such as gene loss often have no discernable phenotypic effects,  
299 underscoring the roles of genetic redundancy and context-dependence in phenotype generation.  
300 Here, I explored the extent of potential gene loss in *Caenorhabditis* nematodes by describing  
301 orthologous genes that are present (or detectable) in all species but one at a given phylogenetic  
302 level. I then situated these genes within their functional context by connecting them to known  
303 phenotypic roles through *C. elegans* WormBase ontologies and transcriptional data (Boeck, et al.  
304 2016). What can these patterns reveal about the evolution of developmental systems?

305 *Caenorhabditis* species are notable for their morphological constancy in the face of tremendous  
306 genetic divergence (Kiontke, et al. 2004) (although the fig-associated *C. inopinata* is  
307 morphologically exceptional (Kanzaki, et al. 2018; Woodruff, et al. 2018)). Within the *Elegans*  
308 group, species are largely morphologically indistinguishable (although there are male tail  
309 features that distinguish some clades), and mating tests or molecular barcoding is usually  
310 necessary to delineate groups (Kiontke, et al. 2011; Félix, et al. 2014; Stevens, et al. 2019). This  
311 morphological conservation also holds throughout development—the pattern of cell divisions  
312 from the single-cell zygote to the reproductive adult is also largely invariant among species  
313 (Zhao, et al. 2008; Memar, et al. 2018). This conserved developmental pattern persists despite  
314 high genetic divergence; *Caenorhabditis* species span a genetic distance comparable to that  
315 between mice and lampreys (Kiontke, et al. 2004). Here, we document widespread species-  
316 specific gene loss, with dozens of genes being lost on average per species (Fig. 2-3). These  
317 patterns are consistent with previous observations of sequence-level (Cutter 2008) and gene copy  
318 number (Stevens, et al. 2019) variation among *Caenorhabditis* species, with rampant genetic  
319 turnover underlying a stable developmental and morphological system.

320 Notably, only 36% of gene losses observed in other species are associated with obvious  
321 phenotypes in *C. elegans*, with up to 20% of these (on average) actually being essential for  
322 viability and fecundity. How can genes presumably essential for a conserved developmental  
323 system be lost so often? The phenotype terms used to functionally annotate these orthologs are  
324 derived from the vast background knowledge of *C. elegans* model systems genetics and are  
325 generally informed by mutation or RNAi evidence (Schindelman, et al. 2011). So one possible  
326 explanation for this pattern could be that the developmental genetics of *C. elegans* are  
327 exceptionally idiosyncratic such that the functional annotations derived from this species are not  
328 be widely applicable across the genus. In this case, interpretations regarding the loss of essential  
329 genes would be mistaken because *C. elegans* might just be an unusual species; that is, it is  
330 possible that in most *Caenorhabditis* lineages these orthologous groups are indeed dispensable,  
331 and their essential functions are novel to or derived in *C. elegans*. Several lines of evidence lend  
332 some credibility to this explanation. For one, *C. elegans* is a self-fertile hermaphrodite, a mode  
333 of reproduction that is largely the exception in this group, as only three species in the genus have  
334 hermaphrodites (most are male/female) (Stevens, et al. 2019). These species represent 11%  
335 (3/28) of the assemblies used in this study. As the evolution of selfing has profound  
336 consequences for multiple aspects of an organism's biology (Thomas, et al. 2012), this could  
337 promote an idiosyncratic developmental system not comparable to its close relatives. Further,  
338 most of the evidence used for these phenotype terms are derived from studies using the N2 strain  
339 of *C. elegans*, which is thought to be a laboratory domesticated strain (Sterken, et al. 2015). As  
340 *C. elegans* N2 is biologically exceptional with respect to *C. elegans* as a species (Sterken, et al.  
341 2015), it may not be representative of the genus as a whole. *C. elegans* is also divergent in its  
342 regulation of small RNAs—there is ample variation in susceptibility to RNAi by feeding in

343 *Caenorhabditis* and *C. elegans* is particularly vulnerable (Nuez and Félix 2012). Thus *C. elegans*  
344 may represent an idiosyncratic developmental system whose annotations belie a misleading  
345 interpretation of functional gene loss. Future studies utilizing whole-genome approaches to  
346 genetic function in other *Caenorhabditis* species (Verster, et al. 2014) will help to inform the  
347 extent of functional diversity among orthologous genes in this group.

348 Nevertheless, while it is formally possible that *C. elegans* has a particularly idiosyncratic  
349 biology, a much more plausible explanation for these patterns of gene loss is rampant  
350 developmental system drift. Developmental systems drift occurs when divergent developmental  
351 programs underlie otherwise conserved morphological features (True and Haag 2001). One  
352 reason to suspect *C. elegans* is comparable to its close relatives is that many orthologs do  
353 maintain conserved function across the genus (Haag, et al. 2018), and some genes have deeply  
354 conserved functions through nematode phylogeny (Crook 2014; Haag, et al. 2018; Kasimatis and  
355 Phillips 2018). Thus at least some aspects of the *Caenorhabditis* genetic system are conserved  
356 and underlie a static morphology. Instead, the combination of many genes being lost in a species-  
357 specific manner while being essential for fitness in at least one species is consistent with  
358 widespread, species-specific turnover of genetic function despite morphological stasis.

359 This evolutionary pattern of heterogeneity in essential gene function is consistent with other  
360 more direct functional assays across species. *C. elegans* was among the first metazoans to be  
361 interrogated with genome-wide genetic knockdown via RNAi (Fraser, et al. 2000; Gönczy, et al.  
362 2000; Piano, et al. 2000; Maeda, et al. 2001; Kamath, et al. 2003). Since then, dozens of such  
363 screens have been implemented (E Yanos, et al. 2012), which provides a comparatively  
364 exhaustive picture of genetic function within this species. The application of a similar approach  
365 in a close relative, the hermaphroditic *C. briggsae*, affords an opportunity to test the extent of  
366 functional conservation across orthologous genes directly (Verster, et al. 2014). In this case, over  
367 25% of orthologous genes have divergent functions between the two species, consistent with  
368 widespread functional turnover and developmental system drift (Verster, et al. 2014). This point  
369 is further emphasized by a genome-wide RNAi screen among *C. elegans* wild isolates by (Paaby,  
370 et al. 2015), who found widespread variation in maternal-effect gene knockdown penetrance  
371 suggestive of a developmental system in flux within as well as between species. This is perhaps  
372 best exemplified in the development of the nematode vulva, which has long been the study of  
373 detailed genetic analysis and which displays highly divergent developmental processes and  
374 genetic pathways despite yielding a highly conserved morphological structure across nematode  
375 phylogeny (Haag, et al. 2018). The picture of developmental systems drift that is emerging  
376 within nematodes is consistent with observations from other studies, such as variation in  
377 postzygotic isolating factors among closely-related species (including fruit flies, mammals, birds,  
378 butterflies, monkeyflowers, and other taxa (Coyle and Orr 2004)) and vertebrate limb  
379 development (Shubin and Alberch 1986; Haag and True 2018). Overall, the patterns of gene loss  
380 observed here are consistent with a body of evidence detailing a variety of surprisingly dynamic  
381 developmental processes across the tree of life.

### 382 *Predicting gene loss with transcription and pleiotropy*

383 Functional phenotypic annotations revealed that lost genes often have essential functions in *C.*  
384 *elegans*. Can additional information about these genes be used to predict gene loss? I retrieved  
385 WormBase ontology terms (Lee, et al. 2017), Pfam domains (Finn, et al. 2015), and stage-  
386 specific transcriptional abundance data (Boeck, et al. 2016) to examine if they can differentiate  
387 gene retention from gene loss. With respect to WormBase ontology terms, the number of unique

388 terms per feature for each gene was used and provides crude metrics for the extent of its  
389 pleiotropy (i.e., the number of phenotypes a gene has or the number of tissues and/or  
390 developmental stages a gene is expressed in).

391 Intuitively, one might expect less widely expressed and less pleiotropic genes to be less  
392 constrained by selection and more prone to loss. From a univariate perspective, the impact of  
393 transcriptional abundance and pleiotropy on gene loss varies by phylogenetic scope (Figure 4).  
394 Patterns in genes lost only in *C. inopinata* and in the *Elegans* group largely agreed with these  
395 expectations—retained genes were more likely to be expressed across all developmental stages  
396 (Fig. 4b-4c) and have more WormBase ontology features than lost genes. However, this did not  
397 hold for genes lost when considering the genus as a whole. Surprisingly, genes lost at these  
398 different levels of phylogenetic consideration largely did not overlap (Supplemental Figure 16)  
399 and revealed different patterns of differential transcriptional abundance. Specifically, genes lost  
400 with respect to the *Elegans* group had significant but small effects on transcriptional abundance  
401 across development when compared to retained genes (Figure 4b-c). Conversely, genes lost with  
402 respect to *Caenorhabditis* were not distinguishable from retained genes in transcriptional  
403 patterns (Figure 4b-c). These patterns are largely mirrored in the WormBase term metrics (Figure  
404 4a). As genes lost only in *C. inopinata* exhibited moderately low transcription across the board  
405 (Fig. 4b-c), this suggests that as the phylogenetic scope broadens, the impact of transcription and  
406 pleiotropy (in a single reference species) on gene loss weakens. This is also consistent with  
407 widespread developmental system drift and the rapid evolution of developmental processes.

408 When present, these differences in transcriptional abundance appear to span broad periods of  
409 developmental time, and gene loss at broader phylogenetic levels has miniscule or no effects on  
410 these traits (Fig. 4b-c). In a principal component analysis, retained and loss genes do not overlap  
411 in multidimensional space at any level of phylogenetic scope (Supplemental Figures 26-28).  
412 Furthermore, linear discriminant analysis, whose aim is to find the function that best separates  
413 two groups, is also unable to distinguish lost and retain genes (although prediction is marginally  
414 better in the case of genes lost only in *C. inopinata* (Supplemental Figure 29)). Thus, despite  
415 subtle, broad detectable differences in transcriptional abundance at discrete time points (Fig. 4b-  
416 c; supplemental data), these data cannot predictably distinguish genes with a tendency to be lost  
417 in a species-specific manner, consistent with pervasive turnover of developmental mechanisms  
418 along nematode phylogeny. And although gene loss is difficult to predict in this context, it is  
419 possible that additional biological information (such as those uncovered in the modENCODE  
420 project (Gerstein, et al. 2010)) could be harnessed to understand how and why genes are lost in  
421 this manner.

#### 422 *Caveats*

423 Here, I have set to define and understand patterns of gene loss across *Caenorhabditis* nematode  
424 species with publicly available genome assemblies, and gene losses were inferred through a  
425 common computational pipeline applied to these assemblies and their associated protein sets. A  
426 potential complicating factor in the interpretation of these results is variation in genome  
427 assembly and annotation quality. There is clearly variation in both assembly and annotation  
428 quality in the genomes used in this study (Supplemental Figures 1-5). All genomes in this study  
429 used RNAseq data to inform their annotations (Howe, et al. 2016; Kanzaki, et al. 2018; Yin, et  
430 al. 2018; Stevens, et al. 2019) and provide evidence-based approaches to bolster the reliability of  
431 their protein sets. Despite this, questions remain regarding annotation quality. For instance, *C.*  
432 *sinica* has a notably large protein set with 34,696 coding genes. Inflated gene copy number due

433 to collapsed alleles is a known problem with such hyperdiverse gonochoristic species (Barrière,  
434 et al. 2009), and it is possible that this reflects an overestimate of gene number in this case.  
435 However, overestimates of gene number should not impact inferences of gene loss per se, as  
436 there is no reason to think collapsed alleles would cause biases against annotating genes that are  
437 present. Additionally, BUSCO completeness scores, which are measured by comparing protein  
438 sets against a set of proteins thought to be largely universal among certain organismal groups  
439 (Simão, et al. 2015), reveal the *C. angaria* protein set as an outlier with a 63.5% completeness  
440 score. This is suggestive of an incomplete protein set which would cause overestimates of gene  
441 loss in this species. Thus, particularly for genomics with low completeness metrics, these are  
442 likely to be overestimates of the extent of gene loss in this group.

443 Variation in genome assembly quality is more problematic for this study. Only four of the 28  
444 species used in this study have chromosome level assemblies, and most of the assemblies are  
445 highly fragmented (Supplemental Figures 3-4). Although our computational pipeline can  
446 presumably overcome shortcomings in annotations through genomic alignments, it cannot  
447 account for genomic regions that have not been assembled. Thus a major caveat of this work is  
448 that these specific inferences of gene loss are provisional due to the high variation in  
449 completeness among the genome assemblies used here. Future work using chromosome-level  
450 assemblies will be required to ascertain more precise estimates of gene loss in this group. That  
451 said, these estimates are not without value—most of the assemblies used here have high  
452 completeness metrics (Supplemental Figure 5) and chromosome-level completeness would likely  
453 not have much impact on the qualitative interpretation that essential genes are often lost. This is  
454 consistent with essential genes being lost even in the assemblies with chromosome-level  
455 completeness (Figure 3). So although the quantitative extent of gene loss per species may be  
456 overestimated, the pervasiveness of developmental system drift remains a reasonable  
457 interpretation.

458 Additionally, the method of orthology assignment itself may impose biases upon inferences of  
459 species-specific loss. Here, loss was defined as being present in every species but one, given  
460 some phylogenetic scope. If there is rapid clade-specific genetic divergence, distance-based  
461 clustering may lead to the splitting of orthologous groups. There are thousands of orthologous  
462 groups that are restricted to a few species (see Supplemental Figure 6 for the example of  
463 orthologous groups found only in *C. remanei* and *C. latens*) or are species-specific. Presumably  
464 this can be partly explained by the emergence of clade-specific genes, but this could also be due  
465 to rapid clade-specific divergence. These types of orthologs would be excluded from this  
466 analysis and could actually underestimate the extent of gene loss. Additionally, as the method of  
467 orthologous group inference begins with predicted protein sets, genes that are erroneously  
468 unidentified in multiple species would not be included here, also underestimating the amount of  
469 gene loss. And finally, there is the implicit use of parsimony in assuming gene loss throughout  
470 this study. In all cases gene loss is assumed because orthologs are present in all other species;  
471 there remains the possibility of multiple gene gains for any of these orthologs, although this  
472 parsimony issue is likely not affecting interpretations.

## 473 **Conclusions**

474 *C. elegans* is a widely used model system for biomedical genetics, and it has been at the  
475 forefront of metazoan genomics since the inception of the discipline. However, the organisms  
476 and resources needed to place these findings in their broader evolutionary and phylogenetic  
477 contexts are only recently becoming available. Here, the previously-sequenced genomes of 28

478 *Caenorhabditis* species revealed that all have lost genes that perform essential functions in *C.*  
479 *elegans*, suggesting that developmental processes are rapidly evolving in this group. As  
480 presumably essential genes are turning over rapidly, this also suggests that biological functions  
481 among conserved genes may also be changing quickly. This underscores the need of comparative  
482 approaches in interpreting and translating findings in model systems genetics across large  
483 evolutionary distances.

## 484 **Materials and Methods**

### 485 *Determining species-specific gene loss*

486 28 *Caenorhabditis* protein sets and genome assemblies were retrieved from the *Caenorhabditis*  
487 Genomes Project (Slos, et al. 2017; Stevens, et al. 2019)([caenorhabditis.org](http://caenorhabditis.org); *C. afra*, *C. castelli*,  
488 *C. doughertyi*, *C. inopinata* (formerly *C. sp. 34*), *C. kamaaina*, *C. latens*, *C. monodelphis*, *C.*  
489 *plicata*, *C. parvicauda* (formerly *C. sp. 21*), *C. zanzibari* (formerly *C. sp. 26*), *C. panamensis*  
490 (formerly *C. sp. 28*), *C. becei* (formerly *C. sp. 29*), *C. utelei* (formerly *C. sp. 31*), *C. sulstoni*  
491 (formerly *C. sp. 32*), *C. quickensis* (formerly *C. sp. 38*), *C. waitukubuli* (formerly *C. sp. 39*), *C.*  
492 *tribulationis* (formerly *C. sp. 40*), *C. virilis*) and WormBase Parasite (Howe, et al.  
493 2016)([parasite.wormbase.org](http://parasite.wormbase.org); *C. angaria*, *C. brenneri*, *C. briggsae*, *C. elegans*, *C. japonica*, *C.*  
494 *remanei*, *C. sinica*, *C. tropicalis*), or were otherwise shared (*C. nigoni* (Yin, et al. 2018) and *C.*  
495 *wallacei*, E. Schwarz pers. comm.). The *Diploscapter coronatus* genome (Hiraki, et al. 2017)  
496 was used as an outgroup. See Supplemental Figures 1-5 for measures of assembly size, gene  
497 number, scaffold number, N50, and completeness of these retrieved genome projects.

498 Alternative splice variants were removed from the protein sets such that each protein-coding  
499 gene was represented by the longest-isoform protein. To identify orthologous groups, 841 all v.  
500 all blastp searches (Camacho, et al. 2009) were performed among the protein sets (version  
501 BLAST+ 2.6.0; blastp options -outfmt 6 -evalue 0.001 -num\_threads 8). The blastp results were  
502 then fed into OrthoFinder (Emms and Kelly 2015)(version 1.1.8; options -b -a 10) to define  
503 orthologous proteins. To determine species-specific gene losses, orthologous groups that were  
504 present in every species but one were extracted. It is important to emphasize that throughout this  
505 paper, “loss” will be assumed to be equivalent to this type of species-specific absence, as  
506 opposed to many other possible patterns of repeated loss in multiple species; here we are only  
507 examining patterns of species-specific loss. Additionally, as these orthologous groups were  
508 absent only in one species, loss is the most parsimonious explanation for their absence as  
509 opposed to multiple independent gains in the other species. This analysis was performed at two  
510 phylogenetic levels (all *Caenorhabditis* species and only *Elegans* group species (Figure 1)) as  
511 the number of shared orthologous groups decreases with phylogenetic distance (Supplemental  
512 Figure 6).

513 To be conservative in estimating the extent of species-specific gene loss, additional filters were  
514 applied to the OrthoFinder results. OrthoFinder implements a size-normalization step to BLAST  
515 bit scores in order to account for the correlation between protein length and bit score (Emms and  
516 Kelly 2015). Concerned that poor gene models that inflate gene length would distort the proper  
517 inference of orthologous groups, species-specific orthologous group losses as defined by  
518 OrthoFinder were re-examined for best-reciprocal blastp hits with *C. elegans* among the results  
519 described above. If putative losses were revealed to have a best reciprocal blastp hit with *C.*  
520 *elegans*, they were removed from consideration as such a species-specific gene loss.  
521 Furthermore, as OrthoFinder uses predicted proteins to define orthologous groups, unannotated



522 genes may be spuriously determined as species-specific losses. To address this problem, the *C.*  
523 *elegans* protein constituents of putative losses were aligned to their respective genome  
524 assemblies using tblastn (Camacho, et al. 2009)(version BLAST+ 2.5.1; options -outfmt 6),  
525 which searches entire translated genomes without the need of gene models. If a putative lost *C.*  
526 *elegans* protein had a best tblastn hit outside of a predicted coding gene in the respective genome  
527 assembly (using the BEDtools (Quinlan 2014) *intersect* function (version 2.25.0; option -v) with  
528 the respective genome assembly's annotations to retrieve alignments that fall outside of predicted  
529 protein-coding regions), this ortholog was then removed from consideration as being a species-  
530 specific loss. This pipeline then accounts for problems incurred by poor gene annotations which  
531 OrthoFinder cannot address.

### 532 *Connecting species-specific ortholog losses to WormBase ontology, Pfam domain, and RNAseq* 533 *data*

534 To understand the functional and biological characteristics of lost genes, the *C. elegans* members  
535 of genes lost in all non-*C. elegans* species (at both levels of phylogenetic consideration; Figure  
536 1) were extracted. These were then paired with WormBase (Howe, et al. 2015) ontology  
537 information (specifically phenotype, anatomy, life stage, and reference count), which were  
538 retrieved with the SimpleMine tool (Howe, et al. 2015) with all *C. elegans* protein-coding genes.  
539 WormBase is a database housing information regarding the genetics, genomics, and general  
540 biology of *C. elegans* and other nematodes. Particularly, it has collected from the literature and  
541 scientific community genome-wide, gene-specific, and hand-curated information including: the  
542 biological consequences of mutation and RNAi exposure (“phenotype”); tissue-specific  
543 expression patterns (“anatomy”); temporal expression patterns (“life stage”); genetic and  
544 biochemical interactions with other genes and their encoded proteins (“interaction”); and its  
545 number of scientific papers (“reference count”), among other features (Howe, et al. 2015). These  
546 features have been formalized as genomic ontologies (Lee and Sternberg 2003; Schindelman, et  
547 al. 2011) and were used in this study. Essential genes were defined as any of those with  
548 WormBase phenotypes containing the words “lethal,” “arrest,” “sterile,” or “dead.” This  
549 included 58 unique phenotype terms (see supplemental data for the list of essential phenotypes).

550 Domains were identified in the *C. elegans* protein set using all domains defined by the Pfam  
551 database (Finn, et al. 2015). Seed alignments for the domains were retrieved from the Pfam FTP  
552 site ([ftp://ftp.ebi.ac.uk/pub/databases/Pfam/current\\_release/Pfam-A.seed.gz](ftp://ftp.ebi.ac.uk/pub/databases/Pfam/current_release/Pfam-A.seed.gz)), and hidden markov  
553 models were constructed with HMMER (version 3.1b2; function *hmmbuild*) (Eddy 1998) using  
554 default parameters. Then, the models were used to search for all domains across all *C. elegans*  
555 proteins using HMMER (function *hmmsearch* option -tblout) using default parameters. These  
556 results were then used to determine the number of domains per protein-coding gene in *C.*  
557 *elegans*.

558 In addition, stage-specific quantitative RNAseq data was also paired with the *C. elegans*  
559 members of genes lost in all non-*C. elegans* species. Data from Boeck et al. 2016 (Boeck, et al.  
560 2016), which examined transcript abundance at 30 min intervals of embryonic development in *C.*  
561 *elegans*, was used to capture expression dynamics and the potential for predicting gene loss. This  
562 data set also included expression data for postembryonic larval stages, dauer developmental  
563 stages, males, the soma, and the germ line. Here, the mean depth of coverage per million across  
564 biological replicates was taken as the representative transcript level for a gene at a given stage.  
565 And as only 19,712 protein-coding genes were reported as being expressed in this data set, only  
566 these genes were included in subsequent analyses here.

567 All statistics were performed using the R programming language. Mann-Whitney U tests,  
568 principal components analyses, and MANOVA tests were performed with the base functions  
569 *wilcoxon.test*, *prcomp*, and *manova*. Cohen's *d* effect sizes were estimated using the *cohen.d*  
570 function in the "effsize" package (<https://github.com/mtorchiano/effsize/>). Linear discriminant  
571 analyses were performed using the *lda* function in the "MASS" package (Venables and Ripley  
572 2013), and discriminant functions were projected back onto the data using the *predict* function.  
573 Multiple logistic regression models were performed with the *lrm* function in the "rms" package  
574 <http://biostat.mc.vanderbilt.edu/wiki/Main/Rrms>).

## 575 **Supplementary Material**

576 Supplemental Figure 1. Assembly sizes.

577 Supplemental Figure 2. Number of protein coding genes.

578 Supplemental Figure 3. Scaffold numbers.

579 Supplemental Figure 4. N50 values.

580 Supplemental Figure 5. BUSCO completeness scores.

581 Supplemental Figure 6. The number of recovered shared orthologous groups decreases as more  
582 species are included.

583 Supplemental Figure 7. The distribution of species-specific lost orthologous groups among all  
584 *Caenorhabditis* and only the *Elegans* group.

585 Supplemental Figure 8. BUSCO completeness score and the number of lost orthologous groups  
586 when considering the *Elegans* group.

587 Supplemental Figure 9. BUSCO completeness score and the number of lost orthologous groups  
588 when considering all *Caenorhabditis*.

589 Supplemental Figure 10. Scaffold number and the number of lost orthologous groups when  
590 considering the *Elegans* group.

591 Supplemental Figure 11. Scaffold number and the number of lost orthologous groups when  
592 considering all *Caenorhabditis*.

593 Supplemental Figure 12. N50 and the number of lost orthologous groups when considering the  
594 *Elegans* group.

595 Supplemental Figure 13. N50 and the number of lost orthologous groups when considering the  
596 all *Caenorhabditis*.

597 Supplemental Figure 14. The average gene copy number per orthologous group is low.

598 Supplemental Figure 15. The distribution of orthologous group gene copy numbers less than  
599 five.

600 Supplemental Figure 16. The *C. elegans* gene constituents of species-specific lost orthologous  
601 groups do not largely overlap among different levels of phylogenetic consideration.

602 Supplemental Figure 17. The distribution of number of unique WormBase terms or Pfam  
603 domains among retained and lost genes when considering all *Caenorhabditis*.

604 Supplemental Figure 18. The distribution of number of unique WormBase terms or Pfam  
605 domains among retained and lost genes when considering the *Elegans* group.

606 Supplemental Figure 19. The distribution of number of unique WormBase terms or Pfam  
607 domains among retained and lost genes when considering only *C. inopinata*.

608 Supplemental Figure 20. The distribution of transcriptional abundance by embryonic stage  
609 among retained and lost genes when considering all *Caenorhabditis*.

610 Supplemental Figure 21. The distribution of transcriptional abundance by embryonic stage  
611 among retained and lost genes when considering the *Elegans* group.

612 Supplemental Figure 22. The distribution of transcriptional abundance by embryonic stage  
613 among retained and lost genes when considering only *C. inopinata*.

614 Supplemental Figure 23. The distribution of transcriptional abundance by postembryonic stage  
615 among retained and lost genes when considering all *Caenorhabditis*.

616 Supplemental Figure 24. The distribution of transcriptional abundance by postembryonic stage  
617 among retained and lost genes when considering the *Elegans* group.

618 Supplemental Figure 25. The distribution of transcriptional abundance by postembryonic stage  
619 among retained and lost genes when considering only *C. inopinata*.

620 Supplemental Figure 26. Principal components analysis, all *Caenorhabditis*.

621 Supplemental Figure 27. Principal components analysis, *Elegans* group.

622 Supplemental Figure 28. Principal components analysis, only *C. inopinata*.

623 Supplemental Figure 29. Linear discriminant analysis of transcriptomic and WormBase data.

624 Supplemental Table 1. The top 25 most common phenotypes in WormBase.

625 Supplemental Table 2. The top ten “most pleiotropic” genes as measured by number of unique  
626 WormBase phenotypes.

627 Supplemental Table 3. The top ten most widely studied genes as measured by WormBase  
628 reference count.

629 Supplementary data (essential\_phenotypes.txt; lost\_gene\_list\_all\_caenorhabditis.tsv;  
630 lost\_gene\_list\_elegans\_group.tsv; lost\_gene\_list\_inopinata.txt; lost\_genes\_wormbase\_boeck.tsv;  
631 statistics\_summaries.xlsx;) are available at *Journal*.

### 632 **Data deposition and accessibility**

633 Genome sequences, protein sets, and annotations were retrieved from the *Caenorhabditis*  
634 Genomes Project ([caenorhabditis.org](http://caenorhabditis.org)) or WormBase ParaSite ([parasite.wormbase.org](http://parasite.wormbase.org)). Data files  
635 and code associated with this study have been deposited in Github at  
636 [https://github.com/gcwoodruff/gene\\_loss](https://github.com/gcwoodruff/gene_loss).

### 637 **Acknowledgements**

638 I thank Mark Blaxter, Erich Schwarz, Janna Fierst, Janet Young, Matthew Rockman, and Patrick  
639 Phillips for sharing genomic data. Bill Cresko and his laboratory members provided helpful  
640 comments throughout the development of this work. Patrick Phillips and Erich Schwarz provided  
641 valuable feedback on earlier versions of this manuscript. This work was supported by funding

642 from the National Institutes of health to GCW (5F32GM115209-03) and to Patrick Phillips (R01  
643 GM-102511; R01 AG049396).

## 644 **References**

- 645 Albalat R, Cañestro C. 2016. Evolution by gene loss. *Nature Reviews Genetics* 17:379.
- 646 Altun ZF, Herndon LA, Wolkow CA, Crocker C, Lints R, Hall DH. 2002-2016. WormAtlas.
- 647 Barrière A, Yang S-P, Pekarek E, Thomas CG, Haag ES, Ruvinsky I. 2009. Detecting heterozygosity in  
648 shotgun genome assemblies: Lessons from obligately outcrossing nematodes. *Genome research*.
- 649 Ben-David E, Burga A, Kruglyak L. 2017. A maternal-effect selfish genetic element in *Caenorhabditis*  
650 *elegans*. *Science* 356:1051-1055.
- 651 Bikard D, Patel D, Le Metté C, Giorgi V, Camilleri C, Bennett MJ, Loudet O. 2009. Divergent evolution  
652 of duplicate genes leads to genetic incompatibilities within *A. thaliana*. *Science* 323:623-626.
- 653 Boeck ME, Huynh C, Gevirtzman L, Thompson OA, Wang G, Kasper DM, Reinke V, Hillier LW,  
654 Waterston RH. 2016. The time-resolved transcriptome of *C. elegans*. *Genome research* 26:1441-1450.
- 655 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+:  
656 architecture and applications. *BMC bioinformatics* 10:421.
- 657 Casillas S, Barbadilla A. 2017. Molecular population genetics. *Genetics* 205:1003-1035.
- 658 Chari S, Dworkin I. 2013. The conditional nature of genetic interactions: the consequences of wild-type  
659 backgrounds on mutational interactions in a genome-wide modifier screen. *PLoS genetics* 9:e1003661.
- 660 Consortium DG. 2007. Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* 450:203.
- 661 Corsi AK, Wightman B, Chalfie M. 2015. A Transparent window into biology: A primer on  
662 *Caenorhabditis elegans*. *Genetics* 200:387-407.
- 663 Coyne JA, Orr HA. 2004. *Speciation*: Sinauer.
- 664 Crook M. 2014. The dauer hypothesis and the evolution of parasitism: 20 years on and still going strong.  
665 *International Journal for Parasitology* 44:1-8.
- 666 Cutter AD. 2008. Divergence times in *Caenorhabditis* and *Drosophila* inferred from direct estimates of  
667 the neutral mutation rate. *Molecular biology and evolution* 25:778-786.
- 668 Dietzl G, Chen D, Schnorrer F, Su K-C, Barinova Y, Fellner M, Gasser B, Kinsey K, Oettel S,  
669 Scheiblauer S. 2007. A genome-wide transgenic RNAi library for conditional gene inactivation in  
670 *Drosophila*. *Nature* 448:151.
- 671 Dowell RD, Ryan O, Jansen A, Cheung D, Agarwala S, Danford T, Bernstein DA, Rolfe PA, Heisler LE,  
672 Chin B. 2010. Genotype to phenotype: a complex problem. *Science* 328:469-469.
- 673 E Yanos M, F Bennett C, Kaeberlein M. 2012. Genome-wide RNAi longevity screens in *Caenorhabditis*  
674 *elegans*. *Current Genomics* 13:508-518.
- 675 Eddy SR. 1998. Profile hidden Markov models. *Bioinformatics (Oxford, England)* 14:755-763.
- 676 Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome comparisons  
677 dramatically improves orthogroup inference accuracy. *Genome biology* 16:157.
- 678 Félix M-A, Braendle C, Cutter AD. 2014. A streamlined system for species diagnosis in *Caenorhabditis*  
679 (Nematoda: Rhabditidae) with name designations for 15 distinct biological species. *PLoS One* 9:e94723.
- 680 Ferrari C, Salle R, Callemeyn-Torre N, Jovelín R, Cutter AD, Braendle C. 2017. Ephemeral-habitat  
681 colonization and neotropical species richness of *Caenorhabditis* nematodes. *BMC Ecology* 17:43.
- 682 Fierst JL, Willis JH, Thomas CG, Wang W, Reynolds RM, Ahearne TE, Cutter AD, Phillips PC. 2015.  
683 Reproductive mode and the evolution of genome size and structure in *Caenorhabditis* nematodes. *PLoS*  
684 *genetics* 11:e1005323.
- 685 Finn RD, Coghill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, Potter SC, Punta M, Qureshi M,  
686 Sangrador-Vegas A. 2015. The Pfam protein families database: towards a more sustainable future.  
687 *Nucleic acids research* 44:D279-D285.
- 688 Fraser AG, Kamath RS, Zipperlen P, Martinez-Campos M, Sohrmann M, Ahringer J. 2000. Functional  
689 genomic analysis of *C. elegans* chromosome I by systematic RNA interference. *Nature* 408:325.

690 Gerstein MB, Lu ZJ, Van Nostrand EL, Cheng C, Arshinoff BI, Liu T, Yip KY, Robilotto R, Rechtsteiner  
691 A, Ikegami K. 2010. Integrative analysis of the *Caenorhabditis elegans* genome by the modENCODE  
692 project. *Science* 330:1775-1787.  
693 Gönczy P, Echeverri C, Oegema K, Coulson A, Jones SJ, Copley RR, Dupéron J, Oegema J, Brehm M,  
694 Cassin E. 2000. Functional genomic analysis of cell division in *C. elegans* using RNAi of genes on  
695 chromosome III. *Nature* 408:331.  
696 Good BH, McDonald MJ, Barrick JE, Lenski RE, Desai MM. 2017. The dynamics of molecular evolution  
697 over 60,000 generations. *Nature* 551:45.  
698 Greenberg AJ, Moran JR, Coyne JA, Wu C-I. 2003. Ecological adaptation during incipient speciation  
699 revealed by precise gene replacement. *Science* 302:1754-1757.  
700 Haag ES, Fitch DH, Delattre M. 2018. From “the Worm” to “the Worms” and Back Again: The  
701 Evolutionary Developmental Biology of Nematodes. *Genetics* 210:397-433.  
702 Haag ES, True JR. 2018. Developmental System Drift. *Evolutionary Developmental Biology: A*  
703 *Reference Guide*:1-12.  
704 Hashimshony T, Feder M, Levin M, Hall BK, Yanai I. 2015. Spatiotemporal transcriptomics reveals the  
705 evolutionary history of the endoderm germ layer. *Nature* 519:219.  
706 Hiraki H, Kagoshima H, Kraus C, Schiffer PH, Ueta Y, Kroihner M, Schierenberg E, Kohara Y. 2017.  
707 Genome analysis of *Diploscapter coronatus*: insights into molecular peculiarities of a nematode with  
708 parthenogenetic reproduction. *BMC genomics* 18:478.  
709 Hoballah ME, Gübitz T, Stuurman J, Broger L, Barone M, Mandel T, Dell’Olivo A, Arnold M,  
710 Kuhlemeier C. 2007. Single gene-mediated shift in pollinator attraction in *Petunia*. *The Plant Cell*  
711 19:779-790.  
712 Howe KL, Bolt BJ, Cain S, Chan J, Chen WJ, Davis P, Done J, Down T, Gao S, Grove C. 2015.  
713 WormBase 2016: expanding to enable helminth genomic research. *Nucleic acids research*:gkv1217.  
714 Howe KL, Bolt BJ, Shafie M, Kersey P, Berriman M. 2016. WormBase ParaSite— a comprehensive  
715 resource for helminth genomics. *Molecular and biochemical parasitology*.  
716 Huang W, Massouras A, Inoue Y, Peiffer J, Ràmia M, Tarone A, Turlapati L, Zichner T, Zhu D, Lyman  
717 R. 2014. Natural variation in genome architecture among 205 *Drosophila melanogaster* Genetic Reference  
718 Panel lines. *Genome research*:gr. 171546.171113.  
719 Kamath RS, Fraser AG, Dong Y, Poulin G. 2003. Systematic functional analysis of the *Caenorhabditis*  
720 *elegans* genome using RNAi. *Nature* 421:231.  
721 Kanzaki N, Tsai IJ, Tanaka R, Hunt VL, Tsuyama K, Liu D, Maeda Y, Namai S, Kumagai R, Tracey A,  
722 et al. 2018. Biology and genome of a newly discovered sibling species of *Caenorhabditis elegans*. *Nature*  
723 *communications*.  
724 Kasimatis KR, Phillips PC. 2018. Rapid gene family evolution of a nematode sperm protein despite  
725 sequence hyper-conservation. *G3: Genes, Genomes, Genetics* 8:353-362.  
726 Kiontke K, Gavin NP, Raynes Y, Roehrig C, Piano F, Fitch DH. 2004. *Caenorhabditis* phylogeny predicts  
727 convergence of hermaphroditism and extensive intron loss. *Proceedings of the National Academy of*  
728 *Sciences of the United States of America* 101:9003-9008.  
729 Kiontke KC, Félix M-A, Ailion M, Rockman MV, Braendle C, Pénigault J-B, Fitch DH. 2011. A  
730 phylogeny and molecular barcodes for *Caenorhabditis*, with numerous new species from rotting fruits.  
731 *BMC Evolutionary Biology* 11:339.  
732 Lee RY, Sternberg PW. 2003. Building a cell and anatomy ontology of *Caenorhabditis elegans*.  
733 *Comparative and Functional Genomics* 4:121-126.  
734 Lee RYN, Howe KL, Harris TW, Arnaboldi V, Cain S, Chan J, Chen WJ, Davis P, Gao S, Grove C.  
735 2017. WormBase 2017: molting into a new stage. *Nucleic acids research* 46:D869-D874.  
736 Levin M, Hashimshony T, Wagner F, Yanai I. 2012. Developmental milestones punctuate gene  
737 expression in the *Caenorhabditis* embryo. *Developmental cell* 22:1101-1108.  
738 MacArthur DG, Balasubramanian S, Frankish A, Huang N, Morris J, Walter K, Jostins L, Habegger L,  
739 Pickrell JK, Montgomery SB. 2012. A systematic survey of loss-of-function variants in human protein-  
740 coding genes. *Science* 335:823-828.

- 741 Maeda I, Kohara Y, Yamamoto M, Sugimoto A. 2001. Large-scale analysis of gene function in  
742 *Caenorhabditis elegans* by high-throughput RNAi. *Current Biology* 11:171-176.
- 743 Memar N, Schiemann S, Hennig C, Findeis D, Conradt B, Schnabel R. 2018. Twenty million years of  
744 evolution: The embryogenesis of four *Caenorhabditis* species are indistinguishable despite extensive  
745 genome divergence. *Developmental biology*.
- 746 Nilsson A, Koskiniemi S, Eriksson S, Kugelberg E, Hinton J, Andersson DI. 2005. Bacterial genome size  
747 reduction by experimental evolution. *Proceedings of the National Academy of Sciences* 102:12112-  
748 12116.
- 749 Nowak MA, Boerlijst MC, Cooke J, Smith JM. 1997. Evolution of genetic redundancy. *Nature* 388:167.
- 750 Nuez I, Félix M-A. 2012. Evolution of susceptibility to ingested double-stranded RNAs in *Caenorhabditis*  
751 nematodes. *PLoS One* 7:e29811.
- 752 Paaby AB, White AG, Riccardi DD, Gunsalus KC, Piano F, Rockman MV. 2015. Wild worm  
753 embryogenesis harbors ubiquitous polygenic modifier variation. *Elife* 4:e09178.
- 754 Piano F, Schetter AJ, Mangone M, Stein L, Kempthues KJ. 2000. RNAi analysis of genes expressed in the  
755 ovary of *Caenorhabditis elegans*. *Current Biology* 10:1619-1622.
- 756 Protas ME, Hersey C, Kochanek D, Zhou Y, Wilkens H, Jeffery WR, Zon LI, Borowsky R, Tabin CJ.  
757 2006. Genetic analysis of cavefish reveals molecular convergence in the evolution of albinism. *Nature*  
758 *genetics* 38:107.
- 759 Quinlan AR. 2014. BEDTools: the Swiss - army tool for genome feature analysis. *Current protocols in*  
760 *bioinformatics* 47:11.12. 11-11.12. 34.
- 761 Ren X, Li R, Wei X, Bi Y, Ho VWS, Ding Q, Xu Z, Zhang Z, Hsieh C-L, Young A. 2018. Genomic basis  
762 of recombination suppression in the hybrid between *Caenorhabditis briggsae* and *C. nigoni*. *Nucleic acids*  
763 *research* 46:1295-1307.
- 764 Rödelsperger C. 2018. Comparative genomics of gene loss and gain in *Caenorhabditis* and other  
765 nematodes. In. *Comparative Genomics: Springer*. p. 419-432.
- 766 Schindelman G, Fernandes JS, Bastiani CA, Yook K, Sternberg PW. 2011. Worm Phenotype Ontology:  
767 integrating phenotype data within and beyond the *C. elegans* community. *BMC bioinformatics* 12:32.
- 768 Shimizu KK, SHIMIZU - INATSUGI R, Tsuchimatsu T, Purugganan MD. 2008. Independent origins of  
769 self - compatibility in *Arabidopsis thaliana*. *Molecular ecology* 17:704-714.
- 770 Shubin NH, Alberch P. 1986. A morphogenetic approach to the origin and basic organization of the  
771 tetrapod limb. In. *Evolutionary biology: Springer*. p. 319-387.
- 772 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing  
773 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31:3210-3212.
- 774 Slos D, Sudhaus W, Stevens L, Bert W, Blaxter M. 2017. *Caenorhabditis monodelphis* sp. n.: defining the  
775 stem morphology and genomics of the genus *Caenorhabditis*. *BMC Zoology* 2:4.
- 776 Stankiewicz P, Lupski JR. 2010. Structural variation in the human genome and its role in disease. *Annual*  
777 *Review of Medicine* 61:437-455.
- 778 Sterken MG, Snoek LB, Kammenga JE, Andersen EC. 2015. The laboratory domestication of  
779 *Caenorhabditis elegans*. *TRENDS in Genetics* 31:224-231.
- 780 Stevens L, Félix M-A, Beltran T, Braendle C, Caurcel C, Fausett S, Fitch DH, Frézal L, Kaur T, Kiontke  
781 KC, et al. 2019. Comparative genomics of ten new *Caenorhabditis* species. *Evolution Letters*:1-20.
- 782 Thomas CG, Woodruff GC, Haag ES. 2012. Causes and consequences of the evolution of reproductive  
783 mode in *Caenorhabditis* nematodes. *TRENDS in Genetics* 28:213-220.
- 784 True JR, Haag ES. 2001. Developmental system drift and flexibility in evolutionary trajectories.  
785 *Evolution & development* 3:109-119.
- 786 Venables WN, Ripley BD. 2013. *Modern applied statistics with S-PLUS: Springer Science & Business*  
787 *Media*.
- 788 Verster AJ, Ramani AK, McKay SJ, Fraser AG. 2014. Comparative RNAi screens in *C. elegans* and *C.*  
789 *briggsae* reveal the impact of developmental system drift on gene function. *PLoS genetics* 10:e1004077.

790 Vlad D, Kierzkowski D, Rast MI, Vuolo F, Ioio RD, Galinha C, Gan X, Hajheidari M, Hay A, Smith RS.  
791 2014. Leaf shape evolution through duplication, regulatory diversification, and loss of a homeobox gene.  
792 Science 343:780-783.  
793 Winzeler EA, Shoemaker DD, Astromoff A, Liang H, Anderson K, Andre B, Bangham R, Benito R,  
794 Boeke JD, Bussey H. 1999. Functional characterization of the *S. cerevisiae* genome by gene deletion and  
795 parallel analysis. Science 285:901-906.  
796 Woodruff GC, Phillips PC. 2018. Field studies reveal a close relative of *C. elegans* thrives in the fresh  
797 figs of *Ficus septica* and disperses on its *Ceratosolen* pollinating wasps. BMC Ecology.  
798 Woodruff GC, Willis JH, Phillips PC. 2018. Dramatic evolution of body length due to post-embryonic  
799 changes in cell size in a newly discovered close relative of *C. elegans*. Evolution Letters.  
800 Yin D, Schwarz EM, Thomas CG, Felde RL, Korf IF, Cutter AD, Schartner CM, Ralston EJ, Meyer BJ,  
801 Haag ES. 2018. Rapid genome shrinkage in a self-fertile nematode reveals sperm competition proteins.  
802 Science 359:55-61.  
803 Zhao Z, Boyle TJ, Bao Z, Murray JI, Mericle B, Waterston RH. 2008. Comparative analysis of embryonic  
804 cell lineage between *Caenorhabditis briggsae* and *Caenorhabditis elegans*. Developmental biology  
805 314:93-99.  
806 Zufall RA, Rausher MD. 2004. Genetic changes associated with floral adaptation restrict future  
807 evolutionary potential. Nature 428:847.  
808